

Figure 2: Schematic view of the backpropagation in a neural network. We start by calculating the derivative of the cost \mathcal{L} compared to the output (\hat{y} on the diagram) then we go back in the network by reusing the derivatives calculated previously during the calculation of the new derivatives.
Note: in the vector case, it is necessary to add sums as indicated in the equation (2).

layer can be reused to calculate the gradients with respect to the input and the parameters of this same layer: it suffices to multiply it by a simple local gradient (from the output relative to input or parameters).

The successive calculation of the gradients of the different layers is called the *backward pass*.

Questions

- What seem to be the advantages and disadvantages of the various variants of gradient descent between the classic, mini-batch stochastic and online stochastic versions? Which one seems the most reasonable to use in the general case?
- ★ What is the influence of the *learning rate* η on learning?
- ★ Compare the complexity (depending on the number of layers in the network) of calculating the gradients of the loss with respect to the parameters, using the naive approach and the *backprop* algorithm.
- What criteria must the network architecture meet to allow such an optimization procedure?
- The function SoftMax and the loss of cross-entropy are often used together and their gradient is very simple. Show that the loss can be simplified by:

$$\mathcal{L} = -\sum_i y_i \log(\hat{y}_i) + \log\left(\sum_i e^{\hat{y}_i}\right).$$

- Write the gradient of the *loss* (cross-entropy) relative to the intermediate output \hat{y}

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_i} = \dots \Rightarrow \nabla_{\hat{y}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \hat{y}_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \hat{y}_n} \end{bmatrix} = \dots$$

- Using the *backpropagation*, write the gradient of the loss with respect to the weights of the output layer $\nabla_{W_y} \mathcal{L}$. Note that writing this gradient uses $\nabla_{\hat{y}} \mathcal{L}$. Do the same for $\nabla_{W_h} \mathcal{L}$.

$$\frac{\partial \mathcal{L}}{\partial W_{y,j}} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial W_{y,j}} = \dots \Rightarrow \nabla_{W_y} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial W_{y,1,1}} & \dots & \frac{\partial \mathcal{L}}{\partial W_{y,1,n_h}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathcal{L}}{\partial W_{y,n_y,1}} & \dots & \frac{\partial \mathcal{L}}{\partial W_{y,n_y,n_h}} \end{bmatrix} = \dots$$

- ★ Compute other gradients: $\nabla_{W_h} \mathcal{L}$, $\nabla_{W_k} \mathcal{L}$, $\nabla_{W_x} \mathcal{L}$.

5

$$\begin{aligned} \mathcal{L}(y, \hat{y}) &= -\sum_j y_j \log(\hat{y}_j) - \hat{y}_j \log(-\hat{y}_j) \\ &= -\sum_j y_j \log\left(\frac{e^{\hat{y}_j}}{\sum_k e^{\hat{y}_k}}\right) \\ &= -\sum_j y_j \left[\log(e^{\hat{y}_j}) - \log\left(\sum_k e^{\hat{y}_k}\right) \right] \\ &= -\sum_j y_j \hat{y}_j + \left(\sum_j y_j\right) \log\left(\sum_k e^{\hat{y}_k}\right) \\ &= -\sum_j y_j \hat{y}_j + \log\left(\sum_k e^{\hat{y}_k}\right) \end{aligned}$$

$$\begin{aligned} \mathcal{L} &= -\sum_j y_j \hat{y}_j + \log\left(\sum_k e^{\hat{y}_k}\right) \\ \frac{\partial \mathcal{L}}{\partial \hat{y}_i} &= -y_i + \frac{e^{\hat{y}_i}}{\sum_k e^{\hat{y}_k}} \\ &= -y_i + \hat{y}_i = \hat{y}_i - y_i \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial W_{y,j}} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial W_{y,j}} = \sum_k (\hat{y}_k - y_k) \frac{\partial \hat{y}_k}{\partial W_{y,j}}$$

$$\begin{aligned} \frac{\partial \hat{y}_k}{\partial W_{y,j}} &= \left([W_y]_{kj} \right) = \frac{\partial \hat{y}_k}{\partial W_{y,j}} \left(\sum_{p=1}^{n_h} W_{y,p} h_p \right) \\ &= \begin{cases} h_j & \text{if } i=k \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{y,j}} &= \sum_k (\hat{y}_k - y_k) \frac{\partial \hat{y}_k}{\partial W_{y,j}} \\ &= (\hat{y}_j - y_j) \cdot h_j \end{aligned}$$

$$\nabla_{W_y} \mathcal{L} = \begin{bmatrix} (\hat{y}_1 - y_1) h_1 & (\hat{y}_2 - y_2) h_2 \\ (\hat{y}_j - y_j) h_j & (\hat{y}_n - y_n) h_n \end{bmatrix}$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \hat{y}_j} = \frac{\partial \mathcal{L}}{\partial \hat{y}_j} \cdot 1 = \hat{y}_j - y_j$$

$$\begin{aligned} \frac{\partial \hat{y}_k}{\partial W_{y,j}} &= \frac{\partial (W_{y,j} h_j)}{\partial W_{y,j}} = \frac{\partial (h_j)}{\partial W_{y,j}} = \begin{cases} h_j & \text{if } i=k \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \hat{y}_j - y_j$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \hat{y}_j} = \frac{\partial \mathcal{L}}{\partial \hat{y}_j} \cdot 1 = \hat{y}_j - y_j$$

$$\begin{aligned} \frac{\partial \hat{y}_k}{\partial W_{y,j}} &= \frac{\partial (W_{y,j} h_j)}{\partial W_{y,j}} = \frac{\partial (h_j)}{\partial W_{y,j}} = \begin{cases} h_j & \text{if } i=k \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{\partial \left(\sum_{p=1}^{n_h} W_{y,p} h_p \right)}{\partial W_{y,j}} = h_j \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \sum_k \frac{\partial \mathcal{L}}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \hat{y}_j} = \frac{\partial \mathcal{L}}{\partial \hat{y}_j} \cdot 1 = \hat{y}_j - y_j$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_j} = \sum_k (\hat{y}_k - y_k) \frac{\partial \hat{y}_k}{\partial \hat{y}_j} = (\hat{y}_j - y_j) \cdot 1 = \hat{y}_j - y_j$$

Section 2 - Implementation

We now have all the equations allowing us to make predictions (*forward*), to evaluate (loss) and to learn (*backward* and gradient descent) our model.

We are now going to implement this network with PyTorch. The objective of this part is to gradually familiarize yourself with this framework.

This part is inspired by the PyTorch getting started tutorial available at http://pytorch.org/tutorials/beginner/pytorch_with_examples.html. Do not hesitate to watch it in parallel with this lab to help you write the requested functions.

2.1 Forward and backward manuals

We will start by coding the neural network by simply using basic mathematical operations and therefore directly transcribing our mathematical equations. We will use the functions found in the `torch`: <https://pytorch.org/docs/1.2.0/index.html> package.

In torch, data is stored in objects of type `torch.Tensor`, equivalent to `numpy.array`. The basic PyTorch functions return objects of this type.

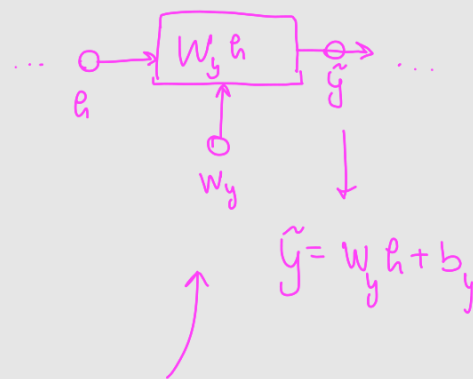
- ★ ★ Discuss and analyze your experiments following the implementation. Provide pertinent figures showing the evolution of the loss; effects of different batch size / learning rate, etc.
- Write the function `init_params(nx, nh, ny)` which initializes the weights of a network from the sizes n_x , n_h and n_y and stores them in a dictionary. All weights will be initialized according to a normal distribution of mean 0 and standard deviation 0.3.
Hint: use the `torch.randn` and `torch.zeros` functions.
- Write the function `forward(params, X)` which calculates the intermediate steps and the output of the network from an input batch X of size $n_{batch} \times n_x$ and weights stored in `params` and store them in a dictionary. We return the dictionary of intermediate steps and the output \hat{Y} of the network.
Hint: we will use `torch.mm` for matrix multiplication, and `torch.tanh`, `torch.exp`, `torch.sum`.
- Write the function `loss_accuracy(Yhat, Y)` which computes the cost function and the precision (rate of good predictions) from an output matrix \hat{Y} (output of `forward`) with respect to a ground truth matrix Y of the same size, and returns the loss \mathcal{L} and the precision `acc`.
Note: We will use the `_, indsY = torch.max(Y, 1)` function which returns the index of the predicted class (or to be predicted) for each example.
Hints: `torch.mean`, `torch.max`, `torch.log`, `torch.sum`.
- Write the function `backward(params, outputs, Y)` which calculates the gradients of the loss with respect to the parameters and stores them in a dictionary.
- Write the function `sgd(params, grads, eta)` which applies a stochastic gradient descent by mini-batch and updates the network parameters from their gradients and the learning step.
- Write the global learning algorithm using these functions.

16)

$$\nabla_{w_y} e = \frac{\partial e}{\partial w_{y,ij}} = \frac{\partial e}{\partial \tilde{y}} \cdot \frac{\partial \tilde{y}}{\partial w_{y,ij}} =$$

VECTOR
CASE

$$= \sum_{k=1}^{n_y} \left(\frac{\partial e}{\partial \tilde{y}_k} \right) \frac{\partial \tilde{y}_k}{\partial w_{y,ij}} =$$

we know,
see answer 15

$$\frac{\partial \tilde{y}_k}{\partial w_{y,ij}} = \frac{\partial ([w_k^y h + b^y])_k}{\partial w_{y,ij}} = \frac{\partial \tilde{y}_k}{\partial w_{y,ij}} \left(\sum_{p=1}^{n_h} w_{kp} h_p \right) =$$

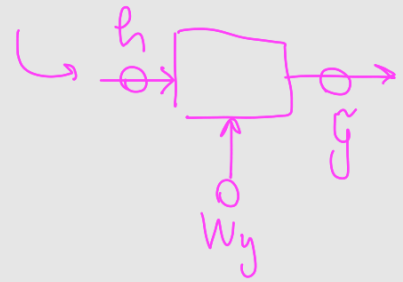
$$= \begin{cases} h_j & \text{if } i=k \\ 0 & \text{otherwise} \end{cases}$$

↓

 $(\hat{y}_i - y_i)$ h_j

$$\frac{\partial e}{\partial w_{y,ij}} = \sum_{k=1}^{n_y} \left(\frac{\partial e}{\partial \tilde{y}_k} \right) \cdot \frac{\partial \tilde{y}_k}{\partial w_{y,ij}} = (\hat{y}_i - y_i) \cdot h_j$$

$$\boxed{\nabla_{b_y} e} = \frac{\partial e}{\partial b_y} = \sum_{k=1}^{n_y} \frac{\partial e}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial b_i}$$



$$\frac{\partial \hat{y}_k}{\partial b_i} = \frac{\partial [w^y \cdot h + b^y]_k}{\partial b_i} = \frac{\partial (b_k)}{\partial b_i}$$

$$= \begin{cases} 1 & \text{if } i=k \\ 0 & \text{otherwise} \end{cases}$$

\Downarrow

$$\frac{\partial e}{\partial b_i} = \frac{\partial e}{\partial \hat{y}_i} = \hat{y}_i - y_i$$

17)

$$\nabla_{\tilde{h}} e$$

$$= \frac{\partial e}{\partial h_i} = \sum_{k=1}^{n_y} \frac{\partial e}{\partial \tilde{y}_k} \cdot \frac{\partial \tilde{y}_k}{\partial h_i}$$

$$\begin{aligned} \frac{\partial \tilde{y}_k}{\partial h_i} &= \frac{\partial ([w^y h + b^y]_k)}{\partial h_i} = \frac{\partial [w^y \cdot h]_k}{\partial h_i} = \\ &= \frac{\partial \left(\sum_{p=1}^{n_h} w^y_{k,p} \cdot h_p \right)}{\partial h_i} = w^y_{k,i} \end{aligned}$$

\Downarrow

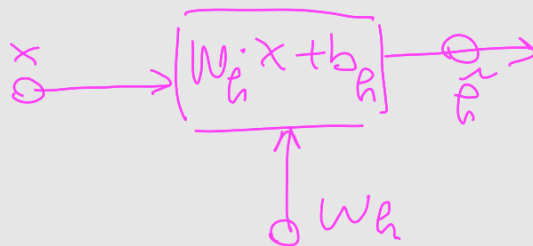
$$\frac{\partial e}{\partial h_i} = \sum_{k=1}^{n_y} \frac{\partial e}{\partial \tilde{y}_k} \cdot \frac{\partial \tilde{y}_k}{\partial h_i} =$$

$$= \sum_{k=1}^{n_y} (\hat{y}_k - y_k) \cdot w^y_{k,i}$$

$$\boxed{\nabla_{w_h} e} = \frac{\partial e}{\partial w_h} = \frac{\partial e}{\partial \tilde{h}} \cdot \frac{\partial \tilde{h}}{\partial w_h} =$$

$$= \sum_{k=1}^{n_h} \left(\frac{\partial e}{\partial \tilde{h}_k} \right) \cdot \left(\frac{\partial \tilde{h}_k}{\partial w_h} \right)^*$$

LOOK PREVIOUS
ANSWER



$$* \frac{\partial \tilde{h}_k}{\partial w_h} = \frac{\partial ([w^h x + b^h]_k)}{\partial w_h} = \begin{cases} x_j & \text{if } i=k \\ 0 & \text{otherwise} \end{cases}$$

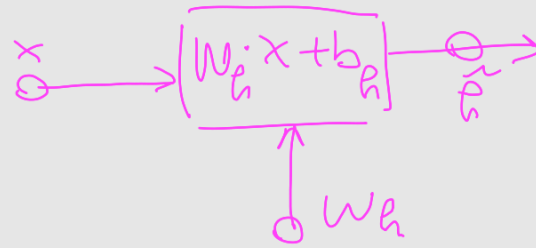


$$\frac{\partial e}{\partial w_h} = \sum_{k=1}^{n_y} (\hat{y}_k - y_k) \cdot w_{k,i} \cdot x_i$$

$$\boxed{\nabla_{b_h} e} = \frac{\partial e}{\partial b_h} = \frac{\partial e}{\partial \tilde{e}^h} \cdot \frac{\partial \tilde{e}^h}{\partial b_h} =$$

$$= \sum_{k=1}^{n_h} \left(\frac{\partial e}{\partial \tilde{e}_k^h} \right) \cdot \left(\frac{\partial \tilde{e}_k^h}{\partial b_h} \right)^*$$

LOOK PREVIOUS
ANSWER



$$* \frac{\partial \tilde{e}_k^h}{\partial b_h} = \frac{\partial ([w^h x + b^h]_k)}{\partial b_h} = \begin{cases} b_i & \text{if } i=k \\ 0 & \text{otherwise} \end{cases}$$

\Downarrow

$$\frac{\partial e}{\partial w_h} = \sum_{k=1}^{n_y} (\hat{y}_k - y_k) \cdot w_{k,i}^y \cdot b_i$$

